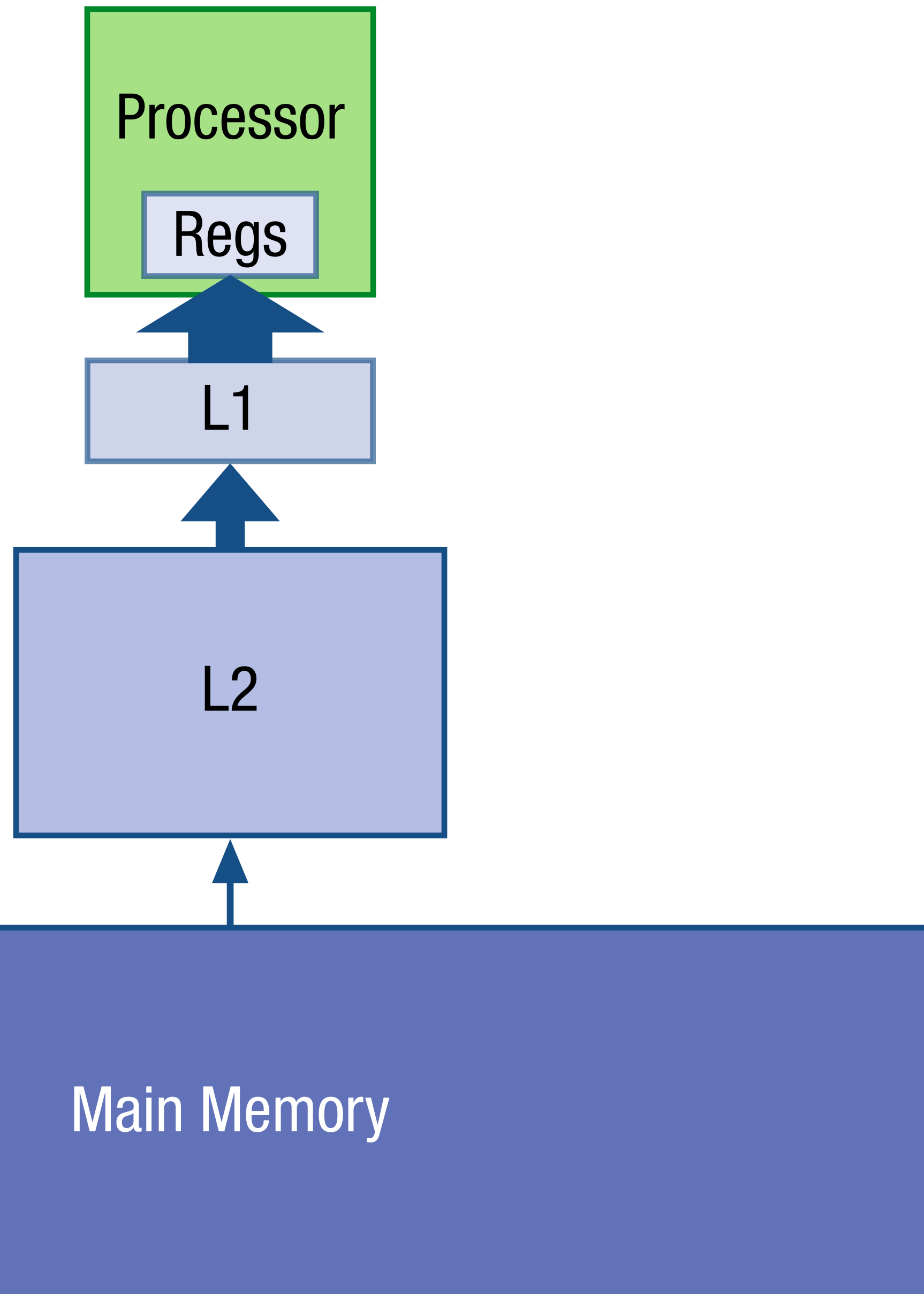# 6.S894
# Accelerated Computing
## Lecture 3: Memory

Jonathan Ragan-Kelley

A Conventional
**Memory Hierarchy**

**Tradeoff:** small, fast, close
**vs.** large, slow, far
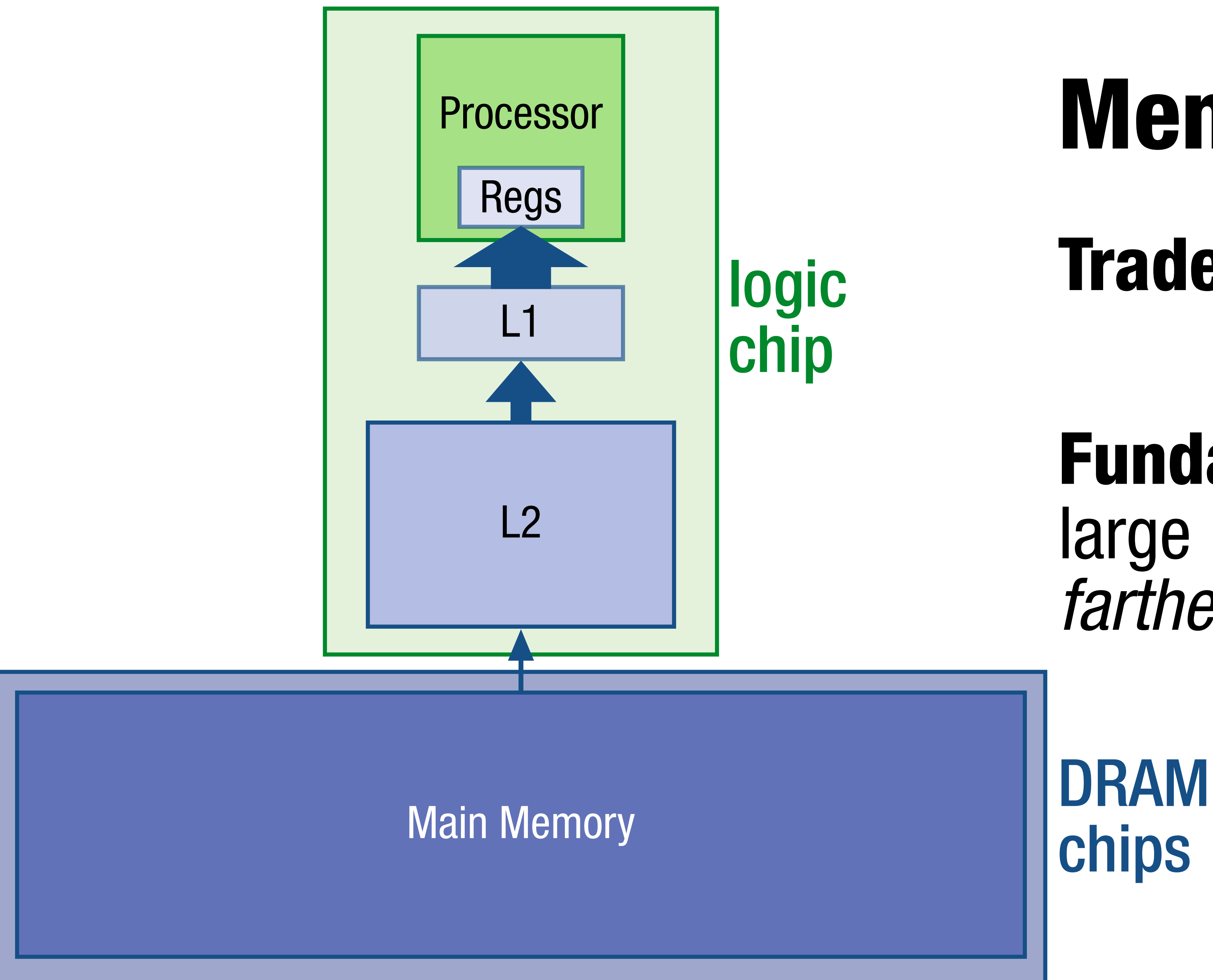
**Fundamental constraint:**
large memories are
*farther away, on average*

Processor

Regs

L1

L2

Main Memory

A Conventional
**Memory Hierarchy**

**Tradeoff:** small, fast, close
**vs.** large, slow, far

**Fundamental constraint:**
large memories are
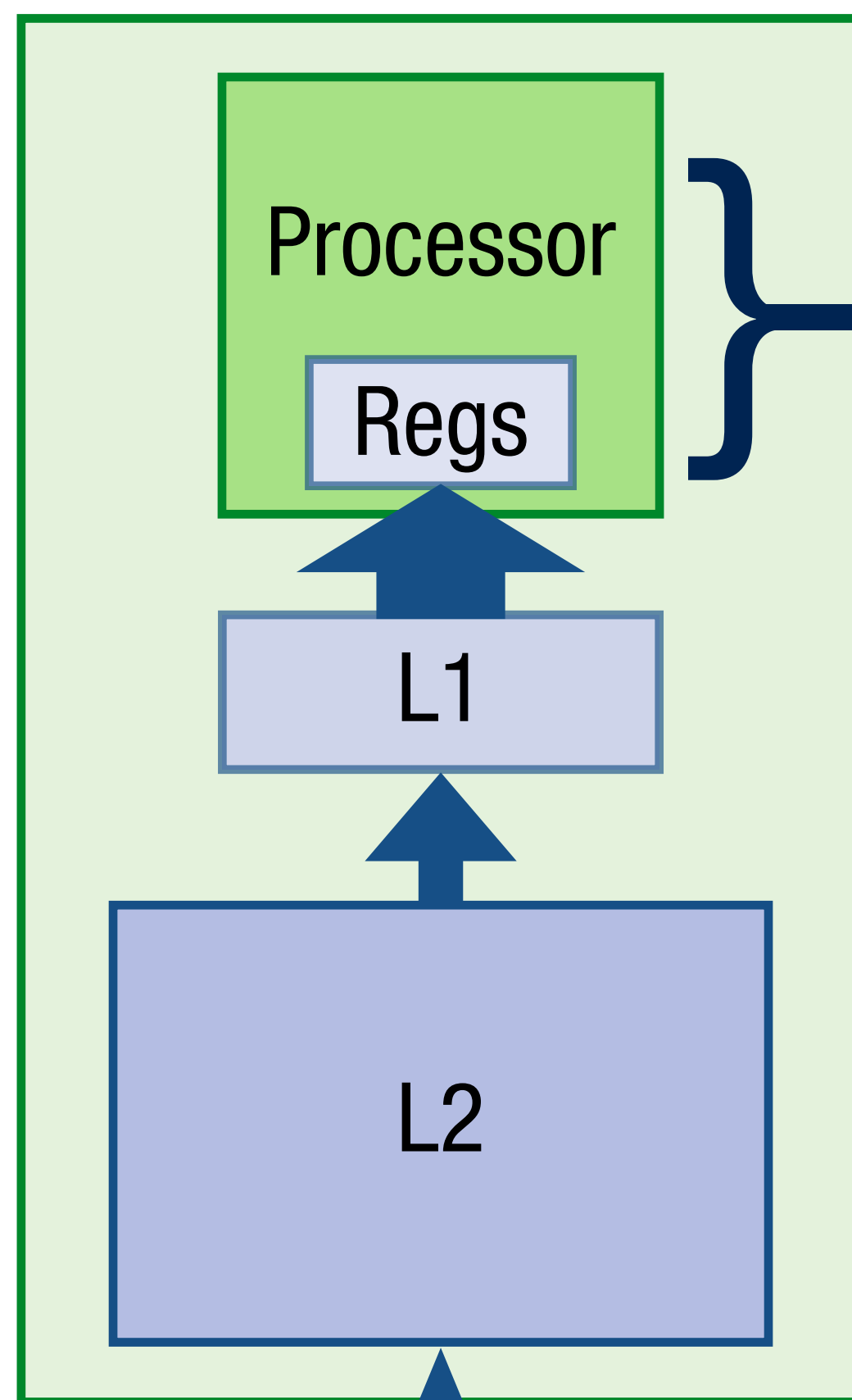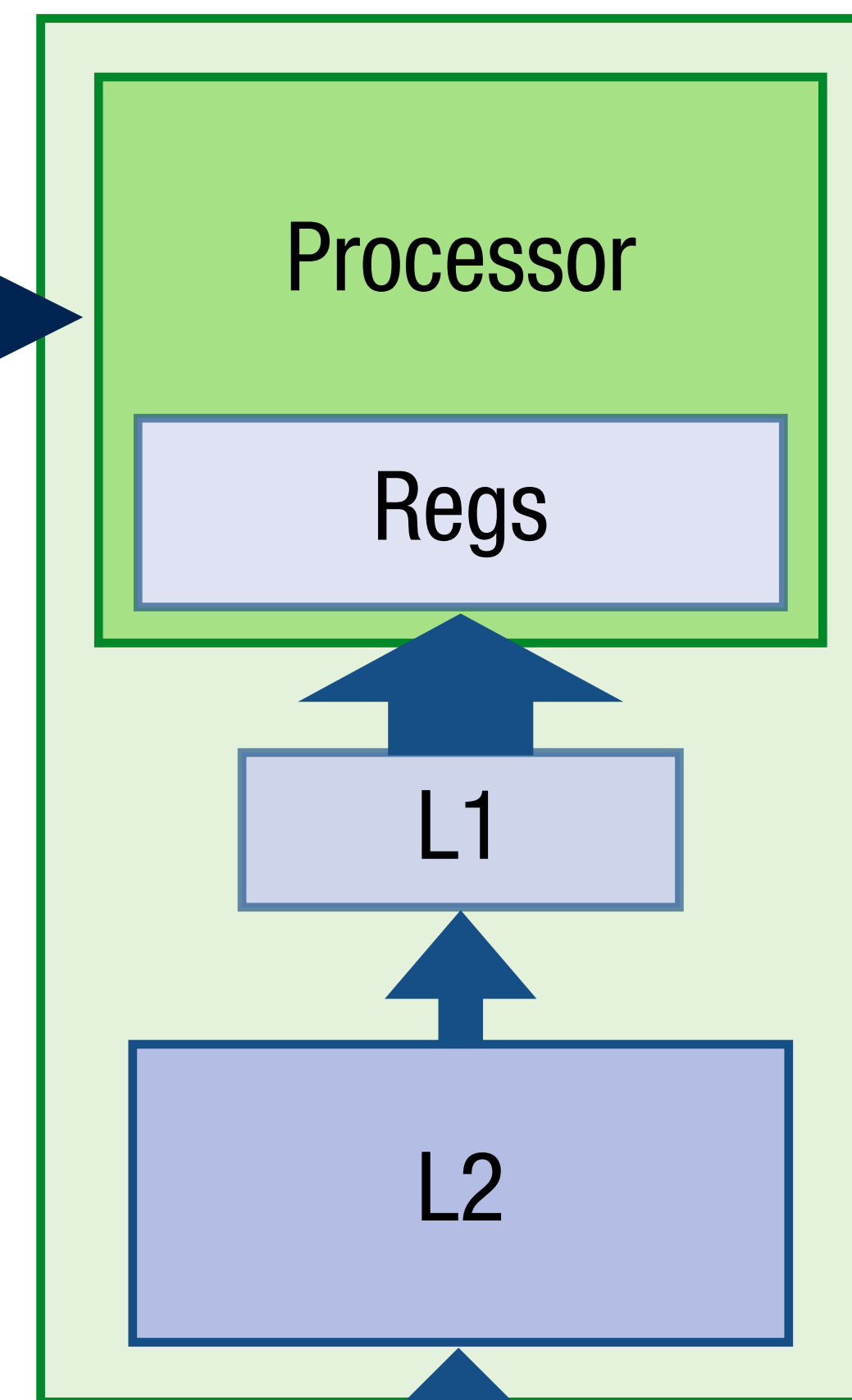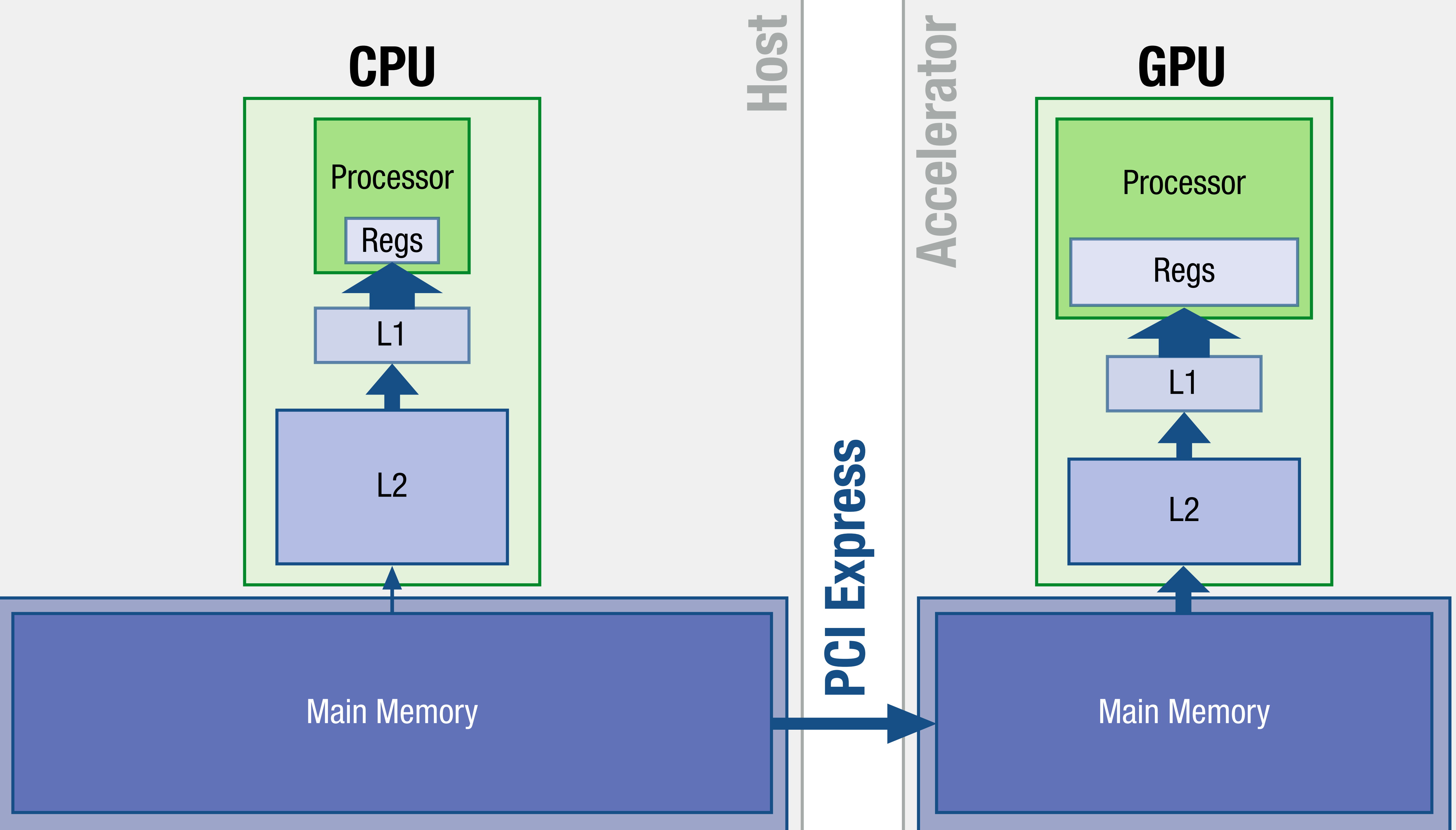*farther away, on average*

**CPU**

Processor
Regs
L1
L2
Main Memory

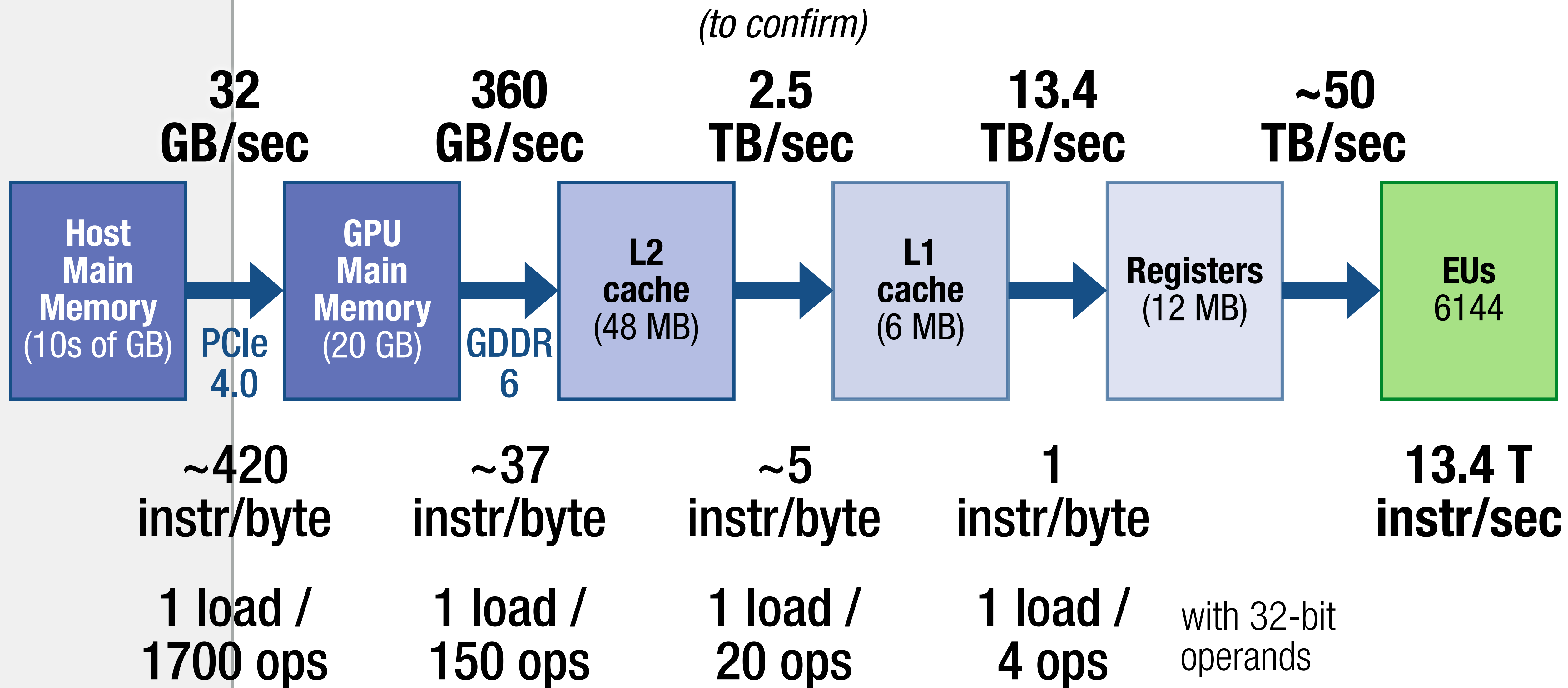more **parallelism**
more **registers**

higher **bandwidth**,
lower **capacity**

**GPU**

Processor
Regs
L1
L2
Main Memory

GPU
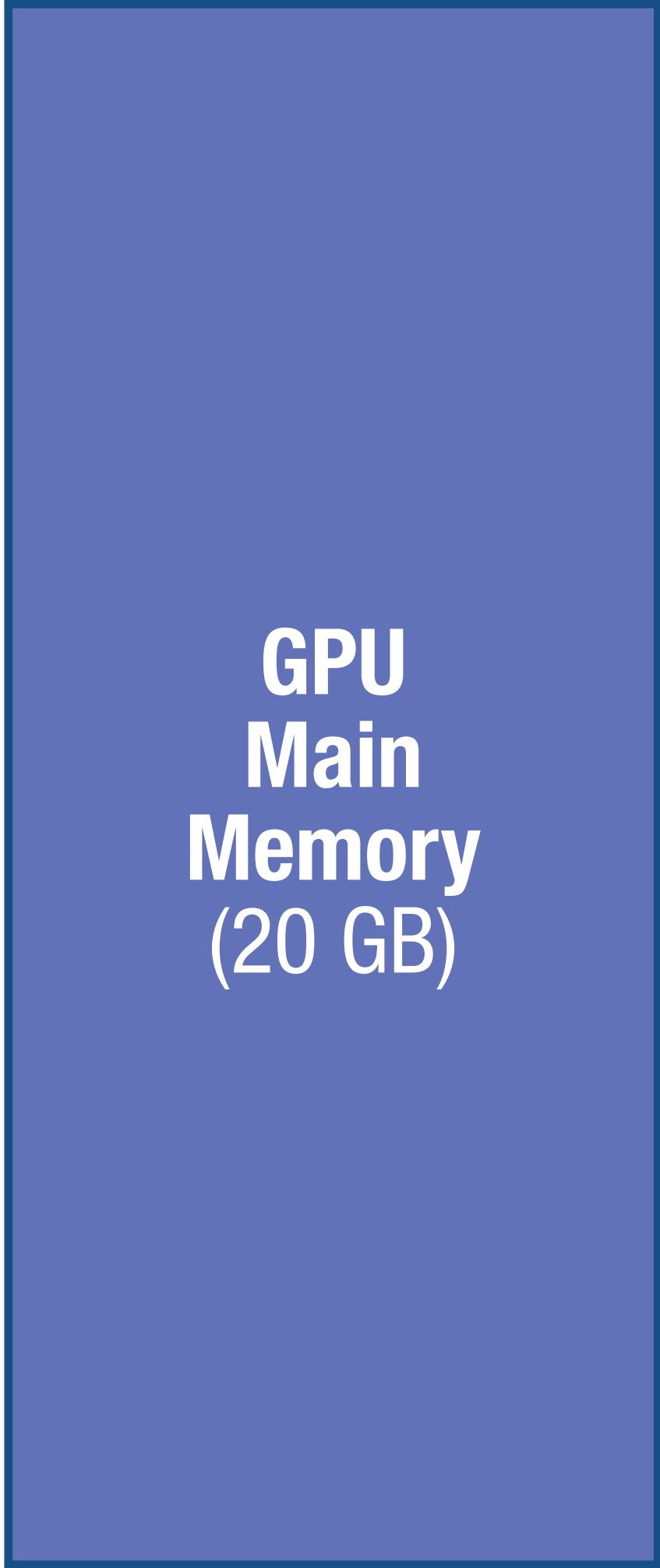Main
Memory
(20 GB)

# How can we get **more bandwidth?**

1. **Faster** bus frequency
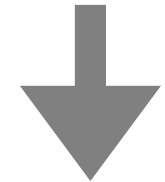
2. **Wider** interface (parallelism)

DRAM
DDR5

**Capacity:** 1-4 GBytes x *n* chips
**Interface:** 32 bits
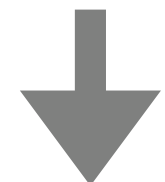**Speed:** 6.4 GT/s

} 25.6 GB/sec

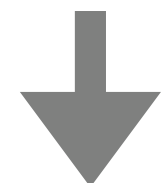(DDR5-6400)

# Increasing memory bus **frequency**

**Tradeoff:**

**faster** signaling

↓

shorter wires

↓

fewer chips

↓

less **capacity**

**DRAM GDDR6**

**Capacity:** 4 GBytes **x _1 chips_**
**Interface:** 32 bits  }
**Speed: 18 GT/s**  } **72 GB/sec**

(GDDR6)

**DRAM DDR5**

**Capacity:** 1-4 GBytes x _n_ chips
**Interface:** 32 bits  }
**Speed: 6.4 GT/s**  } 25.6 GB/sec

(DDR5-6400)
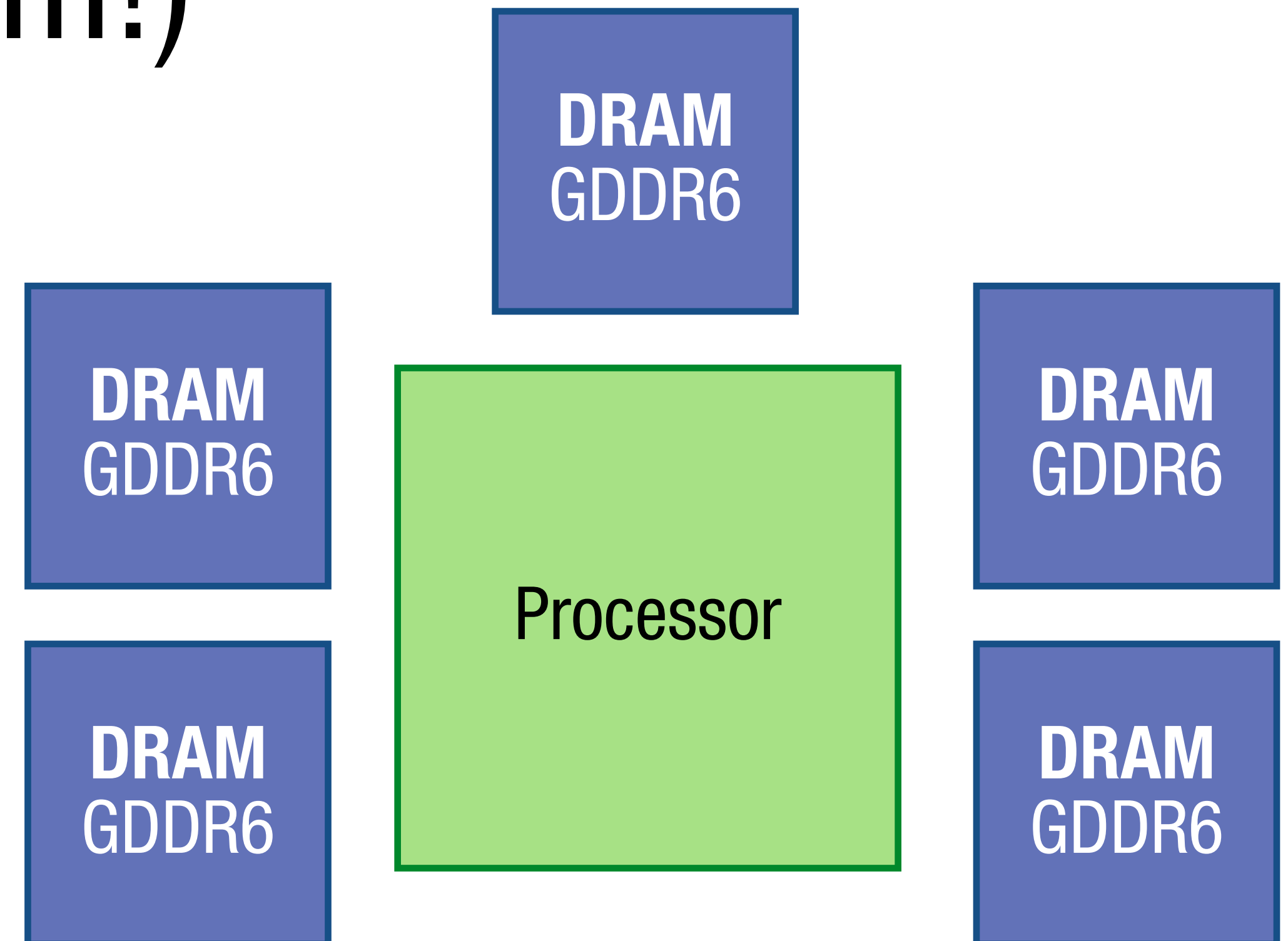
# Increasing memory bus **width** (parallelism!)

**Aggregate:** 5 chips

**Capacity:** 20 GBytes
**Interface:** 160 bits
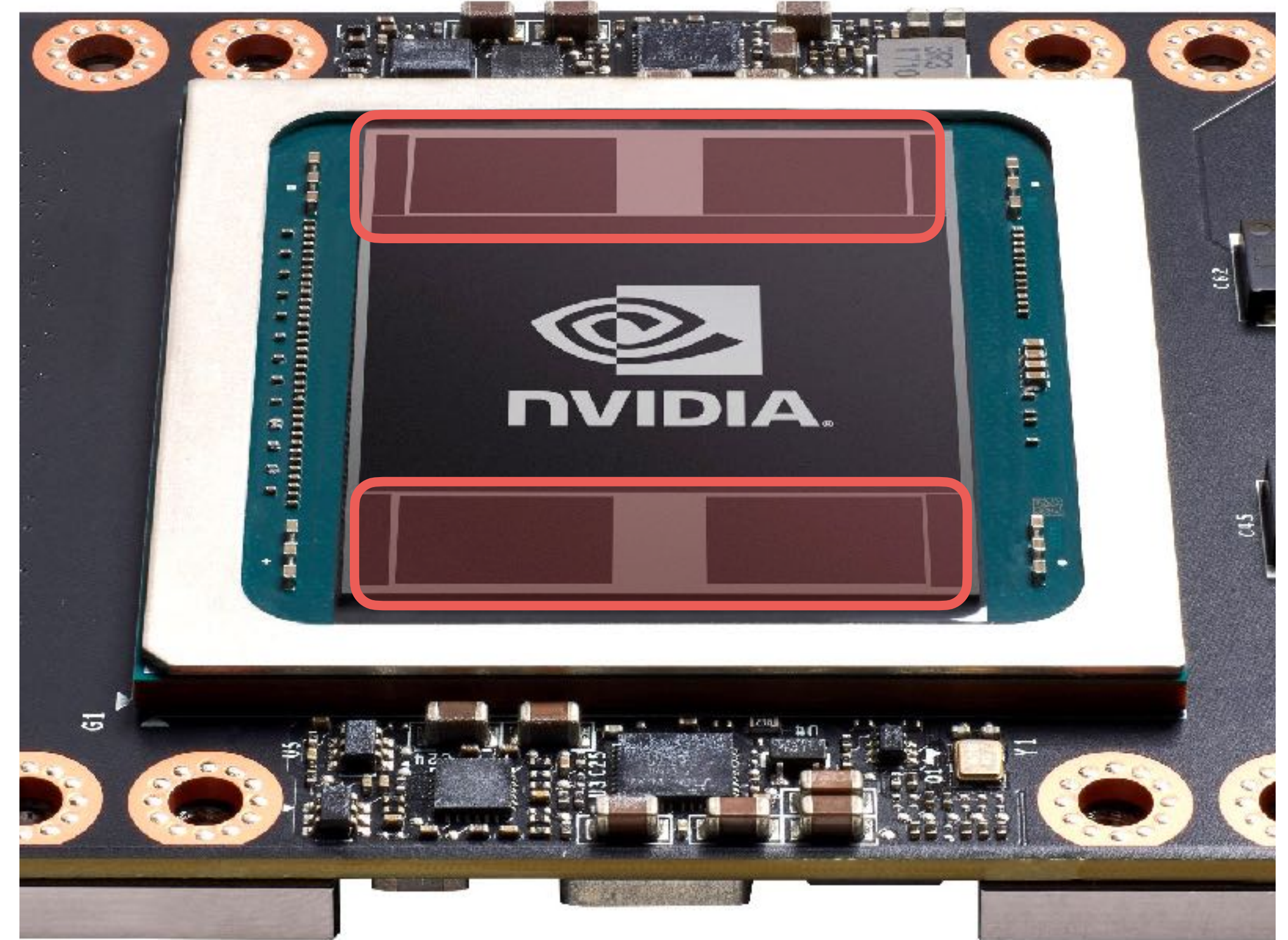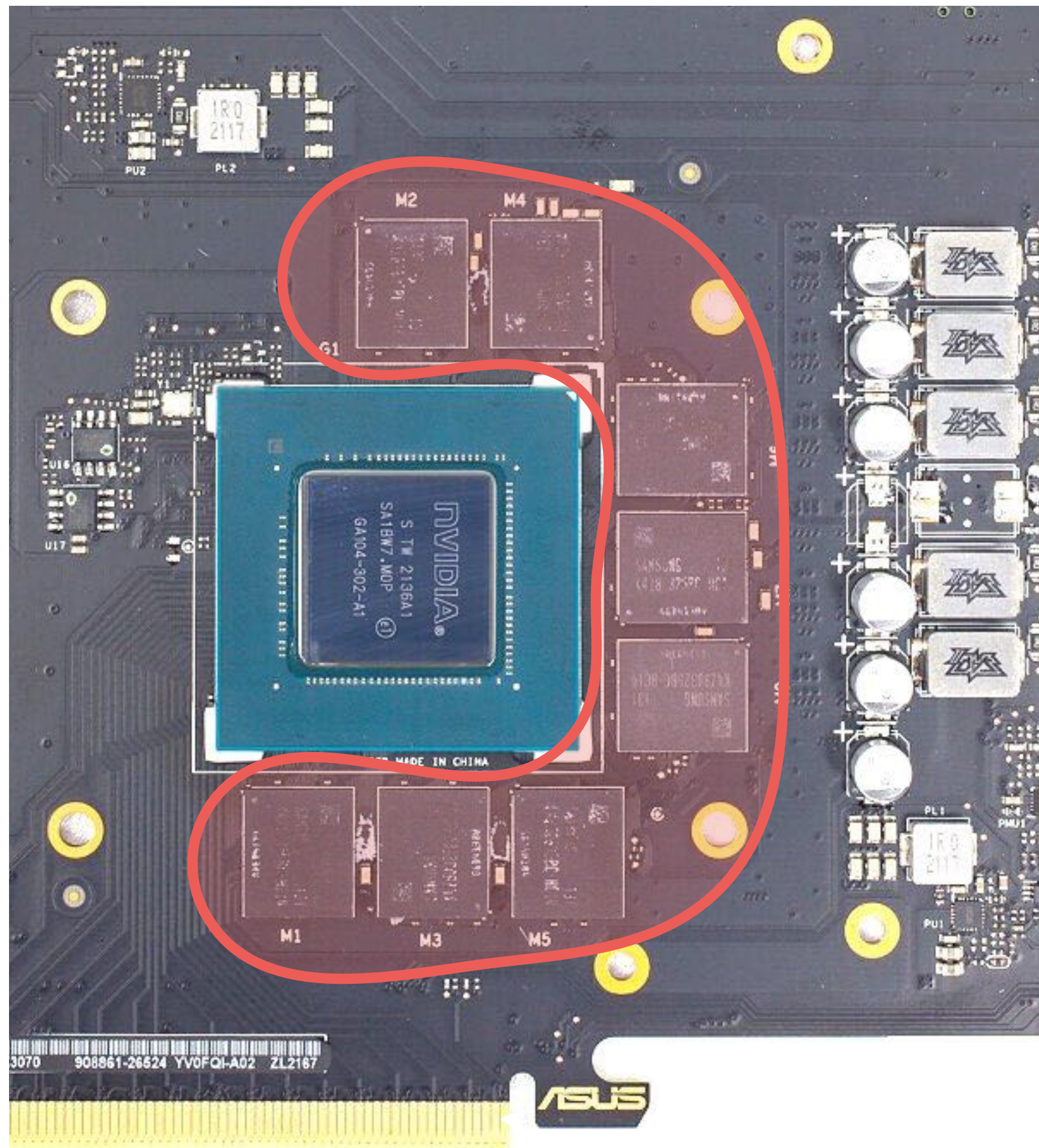**Speed:** 18 GT/s
(GDDR6)
} 360 GB/sec

**Limit:** processor pins
**Practical:** ~384-512 bits
(12-16 chips)

# Increasing memory bus width **further:**
# **package-level** integration

# Increasing memory bus width **further**: **package-level** integration
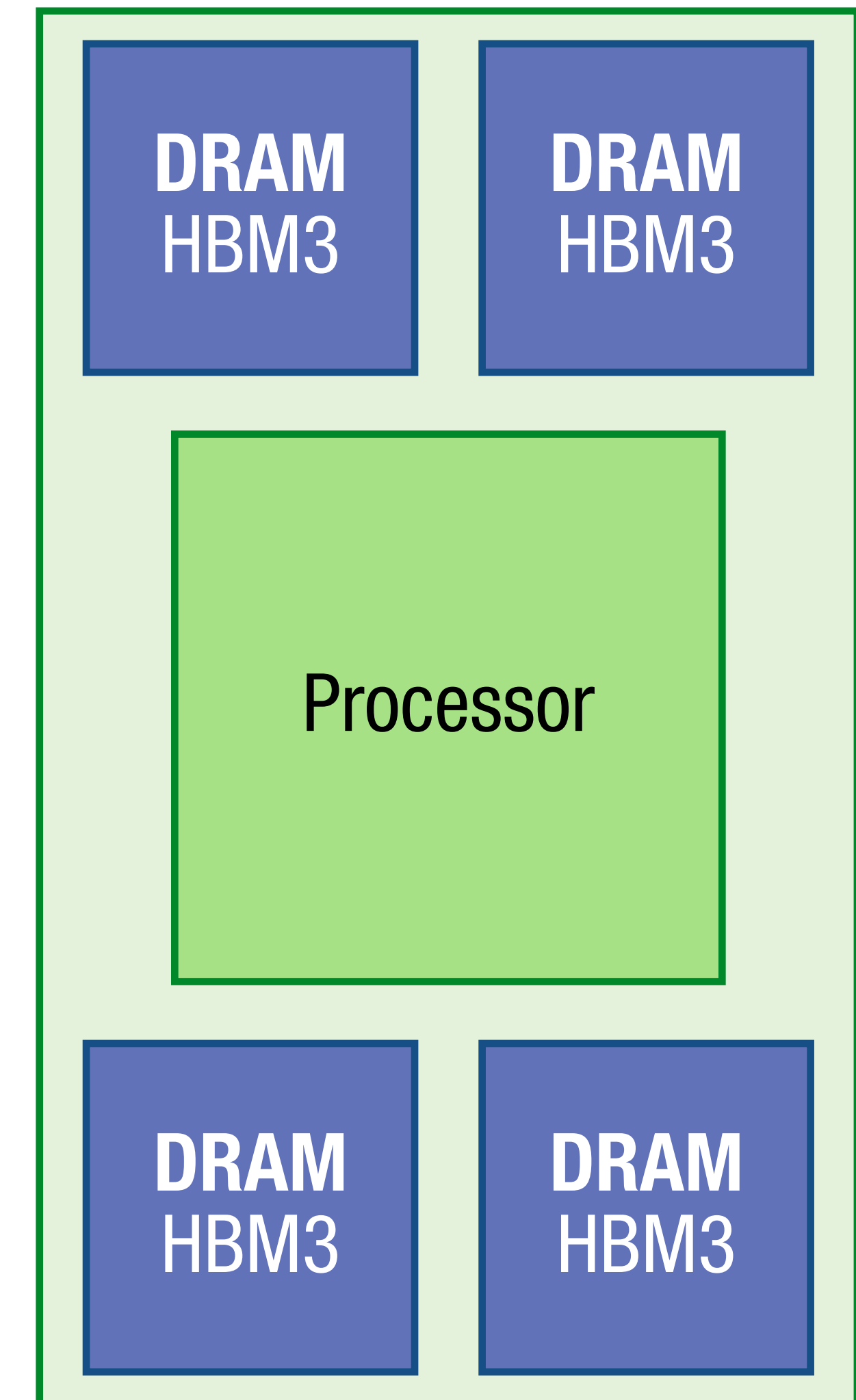
## Per-module:
**Capacity:** 24 GBytes
**Interface:** 2048 bits
**Speed:** 3.2 GT/s } 819 GB/sec
(HBM3)

## Aggregate:
**Capacity:** 96 GBytes
**Interface:** 8192 bits } 3.3 TB/sec

High bandwidth requires
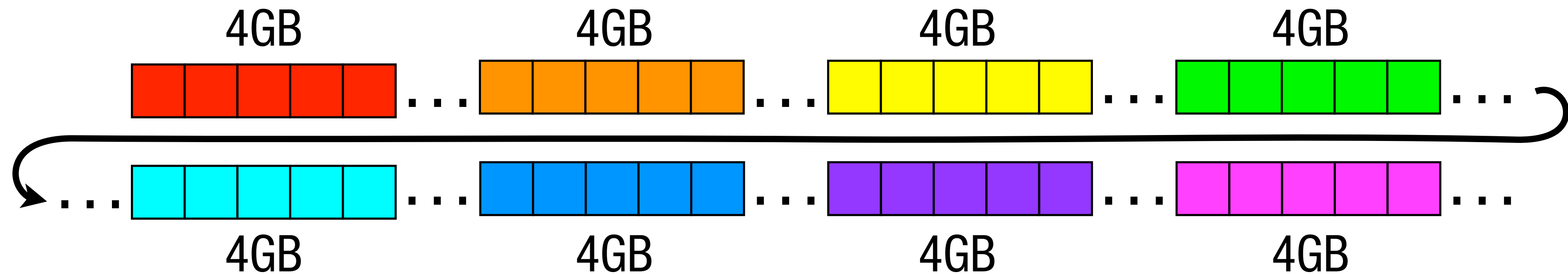**wide memory interfaces**

How can we
**keep them fed?**

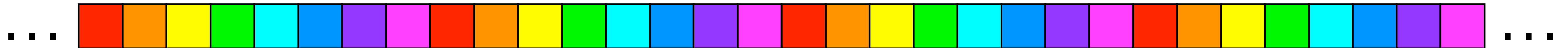# Mapping **address space** to **memory channels**

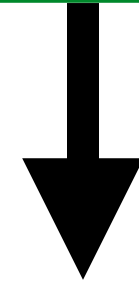# Mapping **address space** to **memory channels**

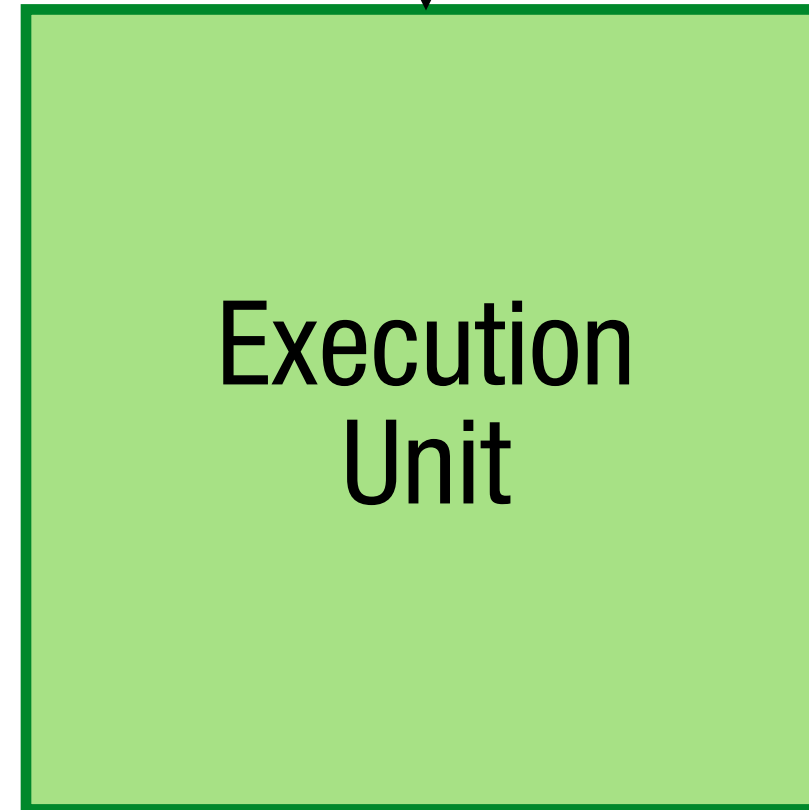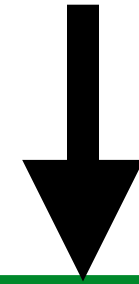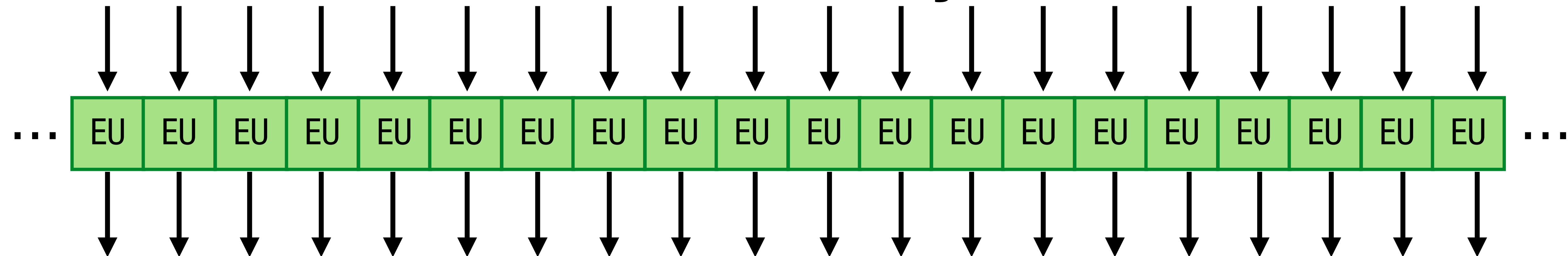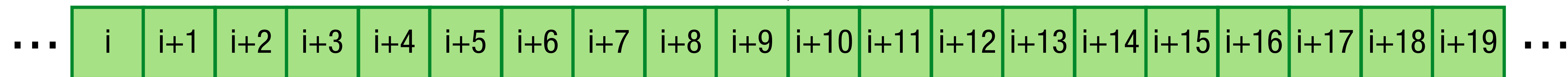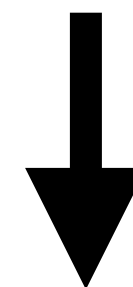# Mapping **address space** to **memory channels**

# 32 bits/cycle

Execution
Unit

# 32 bits/cycle

32 x 32 bits/cycle

...  EU EU EU EU EU EU EU EU EU EU EU EU EU EU EU EU EU EU EU EU  ...

32 x 32 bits/cycle

# Striping memory across channels generally gives **higher bandwidth**